

Texas A&M University - Corpus Christi

Genomics Core Lab

Metabarcoding Manual

Evan Krell, Chris Bird, Martin French, Sharon Magnuson

ABSTRACT

1 INTRODUCTION

The GCL metabarcoding process involves numerous scripts in a variety of programming languages. A metabarcoding project is done by creating a pipeline of these scripts and setting the parameters of each. Because the needs of individual projects differs, it is difficult to create a one-size-fits-all program. Therefore, the use of small, modular programs appears to be the most effective. This manual describes each script in the GCL pipeline, demonstrates an example pipeline, and offers guidance on setting up databases and running on an HPC.

1.1 Charybdis

All of the code resides in a repository called GCL_Charybdis. The repository is located on the HPC at /work/hobi/GCL/GCL_Charybdis. The name resulted from finding a critter from species *Charybdis* from our first set of metabarcoding results.

2 PIPELINE PARAMETERS

You need a number of files and bits of information in order to run the pipeline. The following describes each parameter, and (if relevant), how to obtain them. See section 4 for actual examples of these parameters used in a project.

2.0.1 Project Name (-p)

The name of the project is used as a prefix to many files generated throughout the pipeline. Additionally, the input files need to use this project name in their filenames.

2.0.2 Input Directory (-i)

Certain files need to be placed in this directory before executing the pipeline. Instead of forcing the user to specify several files, they are placed in a single location. The pipeline finds files by using the project name specified with the -p option. The required files for running the generic pipeline script is shown in Table 1.

Table 1. Files required to exist in the input directory.
The variable project name is represented by <projectname>.

<projectname>_forward.fastq	FASTQ with forward reads
<projectname>_reverse.fastq	FASTQ with forward reads
<projectname>_barcodes.txt	Tab-delimited file to match up barcodes to samples

2.0.3 Output Directory (-o)

Files generated by pipeline will be placed in here.

31 **2.0.4 Chunks (-n)**

32 Number of parallel threads to spawn. Use all the cores available if running on an HPC node. But if
33 running on desktop, use less than your number of available cores. For the 40-core workstation in Bird's
34 computer room, I would go with 35.

35 **2.0.5 BLAST Database (-b)**

36 Any valid BLAST database. GCL typically uses the the nucleotide (nt) database from NCBI.

37 `ftp://ftp.ncbi.nlm.nih.gov/blast/db/`

38 **Note that if you use the `-b` option, do not use `-v` option.**

39 **See section 5.4**

40 **2.0.6 BLAST Ignore File (-d)**

41 A text file containing GIs of sequences in the BLAST database to ignore. It is fine to have GIs that are not
42 actually in the BLAST database used. Good for filtering environmental DNA and such.

43 **2.0.7 VSEARCH Database (-v)**

44 VSEARCH is an alternative taxonomic assignment program Rognes et al. (2016). A valid VSEARCH
45 database is a FASTA file. But because FASTA files are a loose standard, not all are compatible. VSEARCH
46 gets unhappy regarding certain formats of the sequence header. An example of a valid FASTA entry is
47 as follows (data is from (Ratnasingham and Hebert, 2007)). Note that the `|` symbol is used to separate
48 metadata field in the header. The first field is a unique sequence ID. The second is a scientific name. The
49 third is the genetic region of the sequence. The fourth is the Barcode of Life Database (BOLD) sequence
50 ID. It could be a NCBI GI or Accession instead.

51 **Note that if you use the `-b` option, do not use `-v` option.**

52 **See section 5.4**

53 `>GBMAA467X14|Pomphorhynchus_laevis|COIX5P|KF559289`
54 `ATGTATGTTTTGGTTGGTGTGTGAGGGGGGCTAATGGGGTTTCTATAAGACTATTAATTCTGA`

55 The script `/work/hobi/GCL/GCL.Charybdis/bin/vsearch_getTAXIDfromBOLDseqid.sh`
56 will get the NCBI taxonomic ID from a BOLD sequence ID.

57 **2.0.8 Chimera Database (-c)**

58 FASTA database used for detection on chimeric reads. Chimera detection can be done using a database or
59 using the *de novo* method Rognes et al. (2016). Currently the pipeline is hard-coded to use the *de novo*
60 method, but still requires a chimera database that will not be used. This issue is described in section 5.3.

61 **2.0.9 Target Sequence Length (-x)**

62 The target sequence length is the length of a read expected using PCR. The forward and reverse primers
63 extract a portion of DNA between them. The value of this option is based on the primers you used, so
64 check their documentation or ask whoever did the PCR.

65 **2.0.10 Charybdis Scripts Directory (-g)**

66 Instead of setting the environment `$PATH` variable or placing the scripts in a systemwide bin directory,
67 the directory holding all the pipeline scripts is an option. Less conventional than `$PATH`-based methods,
68 but very easy to deal with.

69 **2.0.11 NCBI Taxonomy Database (-t)**

70 This database is used to get an organism's scientific name from a numeric taxonomic ID. Taxonomy
71 database downloaded from `ftp://ftp.ncbi.nlm.nih.gov/pub/taxonomy/`.

72 It is the same data as `https://www.ncbi.nlm.nih.gov/Taxonomy/Browser/wwwtax.cgi`

73 **3 PIPELINE ARCHITECTURE**

74 Pipelines are bash scripts that just string together a group of other scripts/programs to perform the metabar-
75 coding process. Example scripts are included in `/work/hobi/GCL/GCL.Charybdis/pipelines`.
76 The script called `charybdis_generic` is expected to handle most projects with little customization.

77 The section will describe all of the steps of the `charybdis_generic` pipeline. Hopefully others
78 will be able to use this information to create a pipeline tailored to other metabarcoding projects.

79 Throughout the text, *attribute* refers to a piece of metadata included in a FASTA sequence header.
80 These are in the form of key-value pairs. Many tools in the OBITOOLS package create and modify
81 attributes, so we made many of our own scripts follow this convention.

82 **3.1 Merge Reads**

83 **Script:** `/work/hobi/GCL/GCL_Charybdis/bin/mergeReads.sh`

84 **Slurm wrapper:** `/work/hobi/GCL/GCL_Charybdis/bin/mergeReads.slurm`

85 The forward and reverse reads are merged into a single FASTA file, and basic quality filtering is done.
86 The reads are matched to sample ID which are included as metadata in the FASTA sequence headers.

87 **3.1.1 Split FASTQ files**

88 In order to run the next sections in parallel, the forward and reverse read FASTQ files are split into equally
89 sized chunks using *fastq-splitter.pl*.

90 **3.1.2 Join forward and reverse reads**

91 After this point, the pipeline uses a single file of merged forward and reverse reads. This is done using
92 OBITOOL's *illuminapairedend* tool. We have it reject reads whose alignment score is under a
93 threshold of 40.

94 **3.1.3 Convert FASTQ to FASTA**

95 Convert the merged files to FASTA. This is done because many tools operate only on FASTA and because
96 we don't need the sequencer quality information after this point. We can always use the unique sequence
97 ID to access it from the origin FASTQ, if needed.

98 **3.1.4 Remove unaligned sequences**

99 We use OBITOOL's *obigrep* to select only those sequences whose alignment mode is "joined". *Obigrep*
100 can search for patterns of values of specified attributes.

101 **3.1.5 Remove sequences whose alignment length is below threshold**

102 We use OBITOOL's *obigrep* to select only those sequences whose alignment length is \geq a threshold. We
103 have this threshold set to 20.

104 **3.1.6 Match sequences to sample ID based on barcode**

105 OBITOOL's *ngsfilter* is used to attach an attribute called *sample* to each sequence. The match is allowed
106 to have 2 sequence errors by setting the flag `-e 2`.

107 **3.1.7 Concatenate into single FASTA**

108 The small chunks are concatenated into a single FASTA file.

109 **3.2 Filter Reads**

110 **Script:** `/work/hobi/GCL/GCL_Charybdis/bin/ObiToolsPipeline.sh`

111 **Slurm wrapper:** `/work/hobi/GCL/GCL_Charybdis/bin/filterReads.slurm`

112 This was the first script written. At the time, the size of the project was not known and this file was
113 intended to be the entire pipeline. Thus, the script is called *ObiToolsPipeline.sh* instead of the
114 more appropriate *filterReads.sh*.

115 **3.2.1 Split into smaller files by sample ID**

116 As before, we form a smaller number of files for parallel processing. However, the sample ID is used to
117 split files instead of an arbitrary size. The sample ID is added to the FASTA filename to keep track.

118 **3.2.2 Remove unneeded attributes**

119 OBITOOLS added a huge number of attributes in a previous step. We get rid of the irrelevant ones with
120 *obiannotate*. Note that attributes are never truly discarded since the sequence ID can always be used
121 to query the intermediate FASTAs generated at each step of the pipeline.

122 **3.2.3 Keep only unique sequences**

123 In order to lower the filesize (often substantially), duplicated sequences are merged into a single represen-
124 tative sequence. The *obimerge* attribute is used to record how many sequences are represented by that
125 sequence. This is done using *obiuniq*.

126 3.2.4 Remove PCR errors

127 PCR errors are filtered using `obiclean`. After numerous experiments, we found that the following
 128 options were performing well for a set of COI sequences from Gulf of Mexico fish. Results from projects
 129 where we reused these settings have appeared to be fine.

Table 2. Settings for `obiclean`

-r	0.5	Controls when less abundant sequences are labeled as variants of a similar more abundant
-d	1	Max number of differences to be considered variants
-H		Keep only a representative of a set of variants.

130 3.2.5 Remove low read length

131 The amplified sequence should be a certain size. Using our COI primers, we expect the size to be 313
 132 base pairs. However, the actual size is often a bit lower or higher. You supply your base pairs length as a
 133 pipeline argument. The allowed range is that size ± 15 .

134 3.2.6 Remove chimeras

135 VSEARCH is used to remove chimeras using the *de novo* method.

136 3.2.7 Concatenate into single FASTA

137 The small chunks are concatenated into a single FASTA file.

138 3.3 Collapse into OTUs

139 **Script:** `/work/hobi/GCL/GCLCharybdis/bin/ClusterOTU.sh`

140 **Slurm wrapper:** `/work/hobi/GCL/GCLCharybdis/bin/ClusterOTU.slurm`

141 3.3.1 Cluster with CROP

142 CROP clusters sequences into OTUs (Hao et al., 2011). CROP has a number of parameters that control
 143 the OTU assignment. Assigning these parameters is tricky, but the CROP manual has a recommended
 144 process for determining the values (TingChenLab, 2017). We initially used the recommendation, but then
 145 modified it to get better results. **However, we don't currently remember how we figured out some of the**
 146 **values used in our process.** You can check the manual and compare to ours, but note that `-b` and `-z` should
 147 only affect performance, not clustering. However, we set `-e` based on `-z` (following the manual), which
 148 does affect results. **We need to investigate further.**

Table 3. Settings for CROP.

NUMSEQS is the total number of sequences. *BP* is the target sequence length.

-r	5	Clusters of size <i>leqr</i> are considered rare
-s		Specifies that 97% similarity needed to cluster
-b	$NUMSEQS/50$	Number of blocks
-z	$(150000/BP) - 0.1 * (150000/BP)$	Max number of sequence a block can hold
-e	$[-z] * 10$	Number of MCMC iterations

149 3.3.2 Correct the OTU size

150 CROP records what sequences end up in an OTU, as well as the total number of sequences represented
 151 by an OTU. However, CROP has no knowledge that `obiclean` and `obimerge` already collapsed
 152 duplicate/similar sequences into. We call `CROP_size_fix.sh` to get the actual OTU counts.

153 3.4 Taxonomic Assignment

154 Taxonomic assignment can be done using either BLAST or VSEARCH. (SAP coming soon). The exact
 155 scripts and slurm wrappers depend on method selected.

156 3.4.1 Run taxonomic assignment program in parallel

157 Split sequences into chunks and BLAST in parallel. Uses the BLAST ignore list (provided in parameters)
 158 to skip assignment to unwanted sequences. For example, we typically try to avoid assignment to
 159 environmental sequences.

160 3.4.2 Concatenate into single FASTA

161 The small chunks are concatenated into a single FASTA file.

162 3.4.3 Convert to Charon format

163 Charon is the name of our own CSV format for assigned sequences. The purpose is that we can write
164 conversion scripts to have assignments from various taxonomic assignment programs into a standard. The
165 Charon file has columns specified by Table 4.

Table 4. Columns for Charon CSV file

Query Sequence ID	Assignment Sequence ID	Scientific Name	NCBI GI	Number of sequences in OTU
-------------------	------------------------	-----------------	---------	----------------------------

166 3.5 OTUs VS Samples

167 Formerly: *Critters VS Tubes*.

168
169 Will explain soon, but basically just a script that makes a column for each sample and a row for each
170 OTU. The cells have a count of the number of sequences of that OTU in that sample. Also has columns
171 of the scientific names for various phylogenetic levels. Also has extra information columns such as the
172 sequence, BLAST scores, etc. No limit on added extra columns to this CSV file.

173 4 EXAMPLE RUN

174 4.1 Prepare Data

175 4.2 Run Pipeline Script

176 The following bash script executes `charybdis_generic.sh`, an implementation of the GCL pipeline.
177 The code segment shown is a single line of bash, using `\` for clarity to break it into one line per option.

```
178 bash /work/hobi/GCL/GCL_charybdis/pipelines/charybdis_generic.sh \  
179     -p Simons-H1.ATTACTCG-TAATCTTA.L001 \  
180     -i /work/hobi/GCL/20180330-Simons/Simons-H1/in \  
181     -o /work/hobi/GCL/20180330-Simons/Simons-H1/out \  
182     -n 20 \  
183     -b /work/hobi/GCL/db/NCBI.BLAST.DBs/nt \  
184     -d /work/hobi/GCL/db/GI_lists/environmental.NCBI.nucl__SORTED.gi \  
185     -c /work/hobi/GCL/db/BOLD.FULL.fasta \  
186     -x 313 \  
187     -g /work/hobi/GCL/GCL_charybdis/bin \  
188     -t /work/hobi/GCL/db/TAXO
```

189 Each option will be explained, and the first 10 lines of each file will be displayed. This way, you can
190 compare the format of your input files to those used for the example. Do not trust the spacing seen here.
191 Tabs may become spaces, etc. See the official file format specifications.

192 **-p** Name of project.

193 We used `Simons-H1.ATTACTCG-TAATCTTA.L001` to match the names of the input FASTQ files. The important
194 thing is that each project should have a unique name for keeping projects separate/organized.

195 **-i** Input directory.

196 This is where the input files should be placed. A slight abuse of terminology, since `<projectname>.samples.txt`
197 is generated by the pipeline, by reading sample IDs from `<projectname>.barcodes.txt`. Note that the
198 FASTQ sequences were shortened to fit on the page. A sequence of five periods (`.....`) represents a removed
199 segment.

```
200 [ekrell@hpcm Simons-H1]$ ls in  
201     Simons-H1.ATTACTCG-TAATCTTA.L001.barcodes.txt  
202     Simons-H1.ATTACTCG-TAATCTTA.L001.forward.fastq  
203     Simons-H1.ATTACTCG-TAATCTTA.L001.reverse.fastq  
204     Simons-H1.ATTACTCG-TAATCTTA.L001.samples.txt  
205  
206 [ekrell@hpcm Simons-H1]$ head in/Simons-H1.ATTACTCG-TAATCTTA.L001.forward.fastq  
207     @MISEQ01:242:000000000-BLT5R:1:1101:15746:1350 1:N:0:ATTACTCGTAATCTTA
```

```

208 ACTTGATAGACCTCGGGGTGGCCGAAGAACCAGAATAAGTGTGTTGGTAAGAATTGGGTCCTC .....
209 +
210 CCCCCFFFFFSGGGGGGGGGHGGGGGHIIHHHHHHHHHHHHHGHHGHHHHHHHHHHHHHHHH.....
211 @MISEQ01:242:000000000-BLT5R:1:1101:16137:1376 1:N:0:ATTACTCGTAATCTTA
212 ATCACGGGAAGTGGATGAACAGTTTATCTCTCCCCTTGCCGGCAACCTGGCCACGCAGGGGC .....
213 +
214 CDDDDDBCCBCFGGGGGGGGGGHIIHHHHHHHHHGGGHIIHHHGGGGGHIIHHHHHHHGGGGGGGGGG .....
215 @MISEQ01:242:000000000-BLT5R:1:1101:15575:1386 1:N:0:ATTACTCGTAATCTTA
216 GTCCGCGGTACTGGATGAACAGTATACCCCCCTTTAGCAGCAGCCATTGCACATGCAGGTGC .....
217
218 [ekrell@hpcm Simons-HI]$ head -n 5 /Simons-HI/ATTACTCG-TAATCTTA.L001.forward.fastq
219 @MISEQ01:242:000000000-BLT5R:1:1101:15746:1350 1:N:0:ATTACTCGTAATCTTA
220 ACTTGATAGACCTCGGGGTGGCCGAAGAACCAGAATAAGTGTGTTGGTAAGAATTGGGTCCTCCCCA .....
221 +
222 CCCCCFFFFFSGGGGGGGGGHGGGGGHIIHHHHHHHHHHHHHGHHGHHHHHHHHHHHHHHHHGGG .....
223 @MISEQ01:242:000000000-BLT5R:1:1101:16137:1376 1:N:0:ATTACTCGTAATCTTA
224 ATCACGGGAAGTGGATGAACAGTTTATCTCTCCCCTTGCCGGCAACCTGGCCACGCAGGGGCATCA .....
225 +
226 CDDDDDBCCBCFGGGGGGGGGGHIIHHHHHHHHHGGGHIIHHHGGGGGHIIHHHHHHHGGGGGGGGGGHH .....
227 @MISEQ01:242:000000000-BLT5R:1:1101:15575:1386 1:N:0:ATTACTCGTAATCTTA
228 GTCCGCGGTACTGGATGAACAGTATACCCCCCTTTAGCAGCAGCCATTGCACATGCAGGTGCATCT .....
229
230 [ekrell@hpcm Simons-HI]$ head -n 5 /Simons-HI/ATTACTCG-TAATCTTA.L001.barcodes.txt -n 5
231 #exp sample tags forward_primer reverse_primer
232 Simons-HI 1901-B ACAGTG:CCGTCC GGWACWGGWTGA..... TANACYTCNGGRTGN.....
233 Simons-HI 1685-C ACAGTG:GTTTCG GGWACWGGWTGA..... TANACYTCNGGRTGN.....
234 Simons-HI 1759-B ACTGAT:ATGTCA GGWACWGGWTGA..... TANACYTCNGGRTGN.....
235 Simons-HI 1774-X ACTGAT:GCCAAT GGWACWGGWTGA..... TANACYTCNGGRTGN.....
236
237 -o Output directory.
238 Files generated by the pipeline, including numerous intermediate file, are placed here.
239 -n Number of parallel instances.
240 The nodes on TAMUCC's HPC cave 20 cores, so 20 were used.
241 -b BLAST Database.
242 Specify path to BLAST database. We are using the nucleotide (nt) database,
243 which is downloaded from ftp://ftp.ncbi.nlm.nih.gov/blast/db/.
244 -d List of GIs for BLAST to ignore This is a list of GIs, where each GI specifies a sequence in the BLAST
245 database to ignore. This file is mandatory, but may be blank (See Section 5.2).
246 We attempt to remove environmental DNA, so our list has GIs of environmental sequences.
247
248 [ekrell@hpcm Simons-HI]$ head -n 5 ...../environmental.NCBI_nucl__SORTED.gi
249 1000014437
250 1000014439
251 1000014441
252 1000014443
253 1000014445
254
255 -c Chimera Database. We use a FASTA file generated from the BOLD database. But, we ended up doing de novo
256 chimera detection. Currently this file is required and just not being used (See section 5.3).
257 -x Target sequence length (basepairs). This number is based on the primers used during PCR.
258 -g Directory of charybdis scripts.
259
260 ls /work/hobi/GCL/GCL_charybdis/bin
261 \ fastq-splitter.pl
262 AddBlast.slurm filterReads.slurm
263 AddVsearch.slurm mergeReads.sh
264 AssignToMarkers.sh mergeReads.slurm
265 AssignToMarkers.slurm ObiToolsPipeline.sh
266 blast_10custom_to_charon_OTU.R ObiToolsPipeline.slurm
267 BlastMeta.sh ObiToolsPipeline_stats.sh
268 BlastMeta.slurm OTU_CVT_addBlast.R
269 ClusterOTU.sh OTUvsTube.slurm

```

```

266 ClusterOTU.slurm                                split_by_marker.R
267 CombineAndCROP.sh                               vsearch_blast6custom_to_charon-BLASTDB.OTU.R
268 crittersVStubes.OTU.R                           vsearch_blast6custom_to_charon.OTU.R
269 CROP_size_fix.sh                                vsearch_getTAXIDfromBOLDseqid.sh
270 CROP.slurm                                       Vsearch.sh
271 determine_marker.R                               Vsearch.slurm
272 determine_marker.sh

273 -t Directory with NCBI taxonomy database Taxonomy database downloaded from ftp://ftp.ncbi.nlm.
274 nih.gov/pub/taxonomy/.
275 -g Directory of charybdis scripts.

276 [ekrell@hpcm pipelines]$ ls /work/hobi/GCL/db/TAXO
277     Accession2Taxid  delnodes.dmp  gc.prt        merged.dmp    nodes.dmp     taxdump.tar.gz
278     citations.dmp    division.dmp  gencode.dmp   names.dmp     readme.txt

```

279 5 AREAS FOR IMPROVEMENT

280 5.1 Support NCBI Accession IDs

281 NCBI is phasing out GIs (sequence IDs), which this pipeline uses at various steps (NCBI, 2016). Accession IDs are
 282 the current standard.

283 5.2 Pipeline should not require BLAST ignore list

284 Currently this is required. If you don't need it, you can supply a blank file. Not very elegant. Would be simple to
 285 make it optional.

286 5.3 Chimera database should not be mandatory

287 The generic pipeline uses *de novo* chimera detection, but still requires the chimera database. A better solution would
 288 be to default to *de novo*, but use the database if supplied.

289 5.4 Taxonomic assignment methods need not be mutual exclusive

290 Currently, you have to specify either a BLAST or VSEARCH database. Will receive a warning if attempting to supply
 291 both options. But, there is no reason to keep the user from specifying both. The pipeline could simply branch at that
 292 point, do assignment with both, and perform remaining pipeline steps with both results. Work would need to be done
 293 to organize filenames to differentiate. The current problem is that all assignment methods generate files with the same
 294 name. The assignment method itself would need to be part of the filename.

295 6 STATISTICAL ASSIGNMENT PACKAGE (SAP)

296 **We have SAP working, but not meshed into the pipeline yet. So it has its own section, at least for now.**

297 SAP (Munch et al., 2008) differs from typical assignment software (BLAST, VSEARCH) in that it is based on
 298 computing a posterior probability that the query sequence is a member of a particular clade. Two major phases are
 299 involved in an SAP taxonomic assignment. First, BLAST is executed to find a set of similar sequences. Since similar
 300 sequences should indicate similar biological properties, the set of accepted BLAST sequences are called the set of
 301 homologues with the query. To strengthen the assumption that the sequences are homologous, there is a limit on
 302 how distant the furthest apart homologues can be. Several parameters are used to control BLAST's sequence search
 303 and also to control how SAP chooses which of those sequences to include in the homologue set. It can and does
 304 happen that a single species has several sequences in the database. However, a particular taxonomic group appears
 305 only once in the phylogenetic tree. Therefore, SAP focuses on diversity and only uses the best match for a given
 306 species. While this is necessary to construct meaningful trees, it also means that the species which only appears once
 307 in a set of BLAST results is given equal weight to a predominate species result. The set of homologues and the query
 308 sequence are used to sample a large number of trees using Markov Chain Monte Carlo. The posterior probability for
 309 any taxonomic group is based on the proportion of trees in which the query sequence was in that group. A simple
 310 case that shows how this can differ from BLAST occurs when several species from the same genus are present in the
 311 homologue set. Where BLAST would simply list the hits in order of sequence ID, SAP determines that it cannot
 312 reliably be determined what species the query belongs to. A top hit may differ little from the hits just below it, and it
 313 is essentially arbitrary to assume that the top hit is the exact species when it scored closely against other species as
 314 well. The SAP result would give a low posterior probability to any particular sequence, but the genus would have a
 315 very high probability.

316 *The above text taken from GCL's work-in-progress manuscript "An Evaluation of Software for Taxonomic*
 317 *Assignment with DNA Metabarcoding"*

318	6.1 Fixing SAP
319	6.2 Running SAP
320	6.3 Evaluating SAP
321	6.3.1 Categorize SAP results
322	6.3.2 Generate SAP report
323	6.4 SAP outstanding issues
324	6.5 Complete SAP example
325	REFERENCES

326 Hao, X., Jiang, R., and Chen, T. (2011). Clustering 16s rRNA for OTU prediction: a method of unsupervised bayesian
327 clustering. *v*, 27:611–8.

328 Munch, K., Boomsma, W., Huelsenbeck, J. P., Willerslev, E., and Nielsen, R. (2008). Statistical assignment of DNA
329 sequences using bayesian phylogenetics. *Systematic biology*, 57(5):750–757.

330 NCBI (2016). Ncbi is phasing out sequence GIS – here’s what you need to know.

331 Ratnasingham, S. and Hebert, P. D. (2007). bold: The Barcode of Life Data System (<http://www.barcodinglife.org>).
332 *Mol. Ecol. Notes*, 7(3):355–364.

333 Rognes, T., Flouri, T., Nichols, B., Quince, C., and Mahe, F. (2016). VSEARCH: a versatile open source tool for
334 metagenomics. *PeerJ*, 4:e2584.

335 TingChenLab (2017). Crop. <https://github.com/tingchenlab/CROP>.